



# Collection of Thai Handwritten Data

Olarik Surinta

Supervisor: dr. Marco Wiering

Promotor: Prof. dr. Lambert Schomaker



# APS Meeting

- Date: Friday 2, March 2012
- Time: 14:00
- Location: VIP room

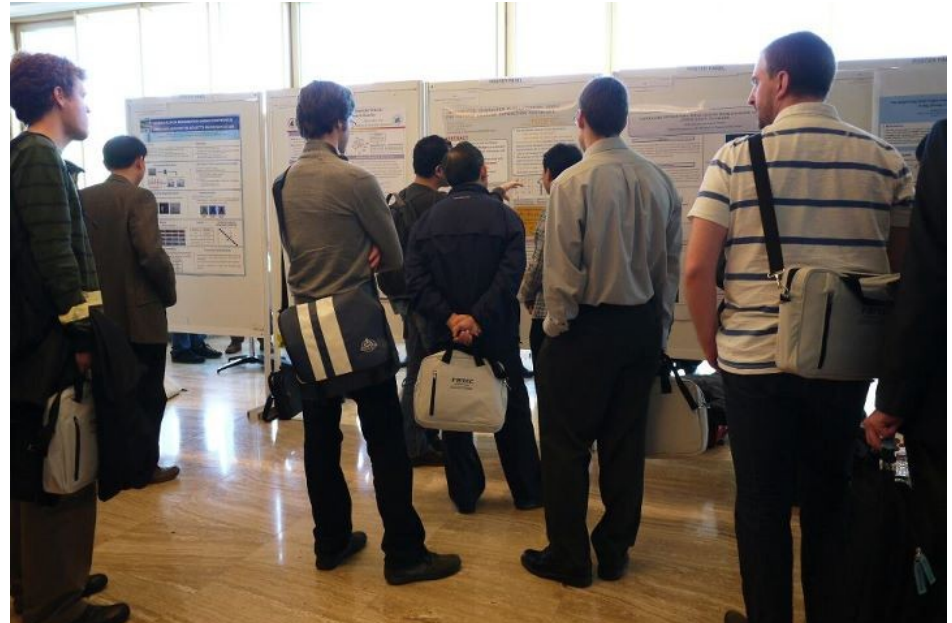
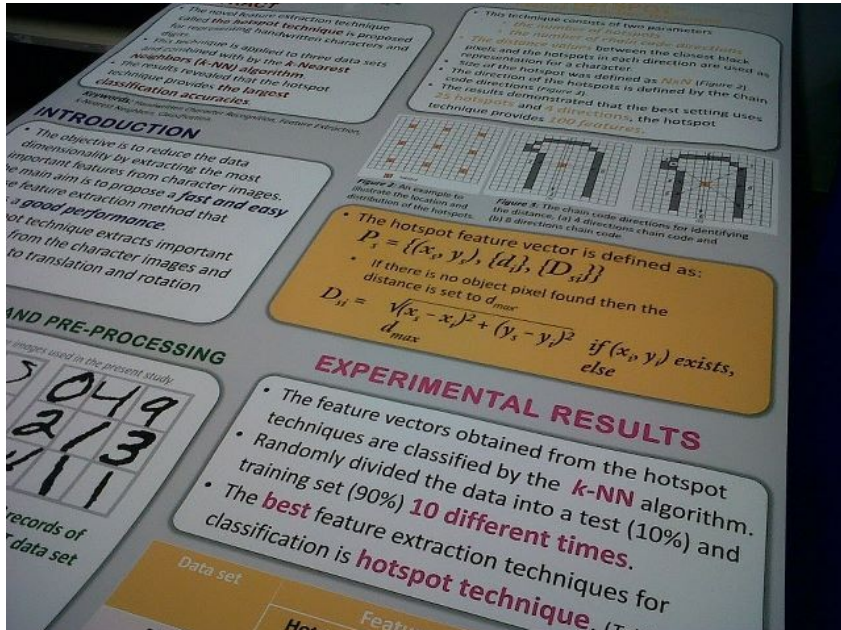


# Previous Event

- 1<sup>st</sup> International Conference on Pattern Recognition Applications and Methods (ICPRAM2012)
- “Handwritten Character Classification using The Hotspot Feature Extraction Technique”



# ICPRAM2012 Event





# Thai Handwritten Dataset

- We used the “Biometrics data for Thai handwritten character's writer form” to collect the handwritten data
- The number of participants are 307 students
- The writers aged from 19-22 years old



Biometrics Data for Thai Handwritten Character's Writer

Page: 1/41

Record Number	Date (DD-MM-YYYY)/time of recording Starting ----- Stopped -----
---------------	--

Sample sentence: (พระบาทสมเด็จพระมงกุฎเกล้าเจ้าอยู่หัว รัชกาลที่ ๖)

Personal contact:

Name	Surname
City of birth	Country of birth
E-mail	Phone number

Personal Cultural:

Nationality
-------------

Personal Biometric:

Gender (female/male)	Age
Handedness (left/right/both)	
Height (centimeter)	Weight (kilogram)

Personal Education:

Mother language	Which's year of study
University	Faculty
Department	

Input Device:

Type of input (pen/pencil/ink pen)	Colour (black/blue/red)
Size of input	

Record Number: \_\_\_\_\_

Page: 2/41

ก	
---	--

ข	
---	--



Biometrics Data for Thai Handwritten Character's Writer

Page: 1/41

Record Number <i>Acc 003</i>	Date (DD-MM-YYYY)/time of recording Starting <i>06 / 12 / 2554</i> Stopped <i>25 / 12 / 2554</i>
---------------------------------	--

Sample sentence: (พระบาทสมเด็จพระมงกุฎเกล้าเจ้าอยู่หัว รัชกาลที่ ๖)

*พระบาทสมเด็จพระมงกุฎเกล้าเจ้าอยู่หัว รัชกาลที่ ๖*

Personal contact:

Name <i>Supansa</i>	Surname <i>Saekram</i>
City of birth <i>Buriram</i>	Country of birth <i>Thailand</i>
E-mail <i>stb 53010916602 @ acc.msu.ac.th</i>	Phone number <i>097-4358602</i>

Personal Cultural:

Nationality <i>Thai</i>
----------------------------

Personal Biometric:

Gender (female/male) <i>female</i>	Age <i>31</i>
Handedness (left/right/both) <i>right</i>	
Height (centimeter) <i>155</i>	Weight (kilogram) <i>45</i>

Personal Education:

Mother language <i>Thai</i>	Which's year of study <i>2 / 2554</i>
University <i>Maharakham University</i>	Faculty <i>Faculty of Accountancy and Management</i>
Department <i>Business Information Technology</i>	

Input Device:

Type of input (pen/pencil/ink pen) <i>pen</i>	Colour (black/blue/red) <i>blue</i>
Size of input <i>0.7</i>	

Record Number: *Acc03*

Page: 2/41

ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก
ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก
ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก
ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก
ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก
ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก
ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก
ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก
ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก
ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก	ก

ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข
ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข
ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข
ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข
ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข
ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข
ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข
ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข
ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข
ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข	ข



# Thai Handwritten Dataset

- This dataset is divided into 2 subset,
  - Thai handwritten character dataset
    - 68 classes, 3,131,400 character pictures
  - Thai handwritten numeric dataset
    - 10 classes, 460,500 numeric pictures

# Handwritten Documents

- The huge Thai handwritten documents
- This picture is half of the whole documents
- All of the documents are 307 folds



# Type of Scanner

- **hp scanjet 5590**  
digital flatbed scanner
- **EPSON GT-1500**  
Document Scanner



# Numeric and Character

0 1 2  
3 4 5  
6 7 8  
9

จ ฉ ช ซ ฌ ญ  
ด ต ถ ท ธ น  
ภ ม ย ร ล ว

# MNIST benchmark dataset

- MNIST configurations;
  - *Size*: 28x28 pixels
  - *Training set*: 60,000 records
  - *Test set*: 10,000 records
- We used ***k*-NN** and **SVM** algorithms to compare pixel-based methods with several feature extraction techniques



# Experimental Results

Feature extraction / Algorithm	SVM	k-NN
Pixel-based (Grey Scale)	97.26	97.44
Pixel-based (Black and White Image)	97.13	96.44
Mark Direction	96.01	91.41
Intensity	96.29	95.55
Hotspot	95.29	90.46
Direction of chain code		
-Only Direction	97.6	95.57
-Only Matrix	98.4	61.68
-All Direction Matrix	97.69	94.5



# Next Plan

- Will publish this research “A comparison of Feature extraction techniques and pixel-based techniques for handwritten character classification ” in IEEE Transactions on Neural Networks
- Data collection
  - Thai handwritten focus on word, sentence, and paragraph levels